

March, 1956

BU-64-M

A FACTORIAL ANALYSIS OF VARIANCE MODEL  
AND ITS APPLICATION IN GENETICS

Douglas S. Robson  
Cornell University

Introduction

Kempthorne (1954) has indicated that his construction of a model for genotypic values in a diploid population at equilibrium under random mating is based upon a model for the analysis of variance in factorial experiments. We shall here show the construction of the analysis of variance model and indicate explicitly its application to the genetic problem.

An Analysis of Variance Model

We suppose that individuals in a population are subjected to "treatments" consisting of combinations of levels of  $N$  factors, say factors  $A^1, A^2, \dots$ , and  $A^N$ . The level of factor  $A^i$  assigned to an individual is determined by chance, independently for each  $i = 1, 2, \dots, N$ , and may be either  $A_1^i, A_2^i, \dots$ , or  $A_{m_i}^i$ . The probability that an individual will receive the particular level  $A_{x_i}^i$  of factor  $A^i$  is written

$$(1) \quad P\left\{A_{X_i}^i = A_{x_i}^i\right\} = p_{x_i}^i, \quad \sum_{x_i=1}^{m_i} p_{x_i}^i = 1.$$

We shall follow the conventional practice of using the upper case symbol  $X_i$  to designate the chance variable and the lower case  $x_i$  to denote a specific numerical value which may be taken on by the chance variable. Since we operate only on the subscripts it is useless to carry along the symbol  $A_{x_i}^i$  and we shall write simply  $x_i$ , instead. Thus, we write (1) as

$$P\left\{X_i = x_i\right\} = p_{x_i}$$

for the probability that an individual will receive the level  $x_i$  of the  $i$ 'th factor. By the assumption of independence, the probability that an individual

will receive the particular combination  $x_1, x_2, \dots, x_N$  of levels of the  $N$  factors is

$$(2) \quad P\{X_1=x_1, X_2=x_2, \dots, X_N=x_N\} = P\{X_1=x_1\} P\{X_2=x_2\} \dots P\{X_N=x_N\}.$$

To further simplify the notation we shall denote the set  $(x_1, x_2, \dots, x_N)$  by the symbol  $x$ , so that (2) becomes

$$P\{X = x\} = \prod_{i=1}^N P\{X_i = x_i\}.$$

We now assume the existence of a real- and single-valued function  $g(x)$ , defined for all  $x = (x_1, \dots, x_N)$ , and we shall refer to  $g(x)$  as the average "yield" under "treatment"  $x$ . The construction of an analysis of variance model for the chance variable  $g(X)$  then consists of defining the various "treatment" effects on "yield" in the population. The average effect of a given level  $x_i$  of the  $i$ 'th factor is defined to be the difference between the population mean and the mean of all individuals carrying the level  $x_i$ . Denoting this effect by  $f(x_i)$ , we then have the definition

$$(3) \quad f(x_i) = E\{g(X) | X_i = x_i\} - E\{g(X)\}$$

where the operator  $E$  denotes expectation; thus, by (2),

$$\begin{aligned} f(x_i) &= \sum_{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_N} g(x_1, \dots, x_N) P\{X_1=x_1\} \dots P\{X_{i-1}=x_{i-1}\} P\{X_{i+1}=x_{i+1}\} \dots P\{X_N=x_N\} \\ &\quad - \sum_{x_1, \dots, x_N} g(x_1, \dots, x_N) P\{X_1=x_1\} \dots P\{X_N=x_N\}. \end{aligned}$$

Since  $E\{g(X) | X_i = x_i\}$  depends only on  $x_i$  we shall abbreviate this by writing  $E\{g(X) | X_i = x_i\} = \bar{g}(x_i)$  and, similarly, we define  $E\{g(X)\} = \bar{g}$ , so that (3) becomes

$$(4) \quad f(x_i) = \bar{g}(x_i) - \bar{g}.$$

If the factors  $i$  and  $j$  did not interact in their effect upon yield; i.e., if the effect on yield of any given level  $x_i$  of factor  $i$  were the same regardless of the level of factor  $j$  present in the treatment then the average yield  $\bar{g}(x_i, x_j)$  of all individuals receiving level  $x_i$  of factor  $i$  and  $x_j$  of factor  $j$  would be  $\bar{g} + f(x_i) + f(x_j)$ . The actual difference between  $\bar{g}(x_i, x_j)$  and  $\bar{g} + f(x_i) + f(x_j)$  is therefore defined to be the interaction effect of  $x_i$  and  $x_j$  and is denoted by  $f(x_i, x_j)$ ; thus,

$$\begin{aligned} f(x_i, x_j) &= E\left\{g(X) \mid X_i = x_i, X_j = x_j\right\} - f(x_i) - f(x_j) - \bar{g} \\ (5) \quad &= \bar{g}(x_i, x_j) - f(x_i) - f(x_j) - \bar{g} \end{aligned}$$

or, by (4),

$$(6) \quad f(x_i, x_j) = \bar{g}(x_i, x_j) - \bar{g}(x_i) - \bar{g}(x_j) + \bar{g}.$$

If factors  $i$  and  $j$  interact in their effect upon yield but their interaction effect is the same regardless of the level of factor  $k$  present in the treatment then we would have, for every  $x_k$ ,

$$\begin{aligned} &E\left\{g(X) \mid x_i, x_j, x_k\right\} - [E\left\{g(X) \mid x_i, x_k\right\} - E\left\{g(X) \mid x_k\right\}] - [E\left\{g(X) \mid x_j, x_k\right\} \\ &\quad - E\left\{g(X) \mid x_k\right\}] - E\left\{g(X) \mid x_k\right\} \\ = &E\left\{g(X) \mid x_i, x_j\right\} - [E\left\{g(X) \mid x_i\right\} - E\left\{g(X)\right\}] - [E\left\{g(x) \mid x_j\right\} - E\left\{g(X)\right\}] \\ &\quad - E\left\{g(X)\right\}. \end{aligned}$$

In general, however, the interaction effect of  $x_i$  and  $x_j$  is not constant for all  $x_k$  and the above equality does not hold. We therefore define the interaction effect  $f(x_i, x_j, x_k)$  to be the difference between the left and right sides above; thus,

$$\begin{aligned} f(x_i, x_j, x_k) &= \left( E\left\{g(X) \mid x_i, x_j, x_k\right\} - [E\left\{g(X) \mid x_i, x_k\right\} - E\left\{g(X) \mid x_k\right\}] \right. \\ &\quad \left. - [E\left\{g(X) \mid x_j, x_k\right\} - E\left\{g(x) \mid x_k\right\}] - E\left\{g(X) \mid x_k\right\} \right) \end{aligned}$$

$$\begin{aligned}
 & - \left( E \left\{ g(x) | x_i, x_j \right\} - [ E \left\{ g(x) | x_i \right\} - E \left\{ g(X) \right\} ] \right. \\
 & \quad \left. - [ E \left\{ g(X) | x_j \right\} - E \left\{ g(X) \right\} ] - E \left\{ g(X) \right\} \right) \\
 & = [ \bar{g}(x_i, x_j, x_k) - \{ \bar{g}(x_i, x_k) - \bar{g}(x_k) \} - \{ \bar{g}(x_j, x_k) - \bar{g}(x_k) \} - \bar{g}(x_k) ] \\
 & = [ \bar{g}(x_i, x_j) - \{ \bar{g}(x_i) - \bar{g} \} - \{ \bar{g}(x_j) - \bar{g} \} - \bar{g} ]
 \end{aligned}$$

or

$$\begin{aligned}
 (7) \quad f(x_i, x_j, x_k) &= \bar{g}(x_i, x_j, x_k) - \bar{g}(x_i, x_j) - \bar{g}(x_i, x_k) - \bar{g}(x_j, x_k) + \bar{g}(x_i) \\
 &\quad + \bar{g}(x_j) + \bar{g}(x_k) - \bar{g}
 \end{aligned}$$

or, since  $\bar{g}(x_i, x_j) = f(x_i, x_j) + f(x_i) + f(x_j) + \bar{g}$ , by (5), and  $\bar{g}(x_i) = f(x_i) + \bar{g}$ , by (4), we may write

$$\begin{aligned}
 (8) \quad f(x_i, x_j, x_k) &= \bar{g}(x_i, x_j, x_k) - f(x_i, x_j) - f(x_i, x_k) - f(x_j, x_k) - f(x_i) \\
 &\quad - f(x_j) - f(x_k) - \bar{g}
 \end{aligned}$$

Extending this reasoning, we define the effect  $f(x_{i_1}, \dots, x_{i_v})$  as

$$(9) \quad f(x_{i_1}, \dots, x_{i_v}) = \sum_{t=0}^v (-1)^{v-t} \sum_{(x_{h_1}, \dots, x_{h_t}) \text{ in } (x_{i_1}, \dots, x_{i_v})} \bar{g}(x_{h_1}, \dots, x_{h_t})$$

with the understanding that for  $t = 0$  the set  $(x_{h_1}, \dots, x_{h_t})$  is the (unique) empty set and  $\bar{g}(x_{h_1}, \dots, x_{h_t}) |_{t=0} = \bar{g}$ . To show that (9) may also be written as

$$(10) \quad f(x_{i_1}, \dots, x_{i_v}) = \bar{g}(x_{i_1}, \dots, x_{i_v}) - \sum_{t=0}^{v-1} \sum_{(x_{h_1}, \dots, x_{h_t}) \text{ in } (x_{i_1}, \dots, x_{i_v})} f(x_{h_1}, \dots, x_{h_t})$$

again with the convention that  $f(x_{h_1}, \dots, x_{h_s})|_{s=0} = \bar{g}$ , we observe from (9) that

$$f(x_{h_1}, \dots, x_{h_t}) = \sum_{s=0}^t (-1)^{t-s} \sum_{(x_{k_1}, \dots, x_{k_s}) \text{ in } (x_{h_1}, \dots, x_{h_t})} \bar{g}(x_{k_1}, \dots, x_{k_s})$$

so that

$$\begin{aligned} & \sum_{t=0}^{v-1} \sum_{(x_{h_1}, \dots, x_{h_t}) \text{ in } (x_{i_1}, \dots, x_{i_v})} f(x_{h_1}, \dots, x_{h_t}) \\ &= \sum_{t=0}^{v-1} \sum_{(x_{h_1}, \dots, x_{h_t}) \text{ in } (x_{i_1}, \dots, x_{i_v})} \left\{ \sum_{s=0}^t (-1)^{t-s} \sum_{(x_{k_1}, \dots, x_{k_s}) \text{ in } (x_{h_1}, \dots, x_{h_t})} \bar{g}(x_{k_1}, \dots, x_{k_s}) \right\}. \end{aligned}$$

But every set  $(x_{k_1}, \dots, x_{k_s})$  is a subset of  $\binom{v-s}{t-s}$  different subsets of  $t$  elements from the set  $(x_{i_1}, \dots, x_{i_v})$ , hence the right side above may be written

$$= \sum_{s=0}^{v-1} \sum_{(x_{k_1}, \dots, x_{k_s}) \text{ in } (x_{i_1}, \dots, x_{i_v})} \bar{g}(x_{k_1}, \dots, x_{k_s}) \left\{ \sum_{t=s}^{v-1} (-1)^{t-s} \binom{v-s}{t-s} \right\}$$

and since

$$\sum_{t=s}^{v-1} (-1)^{t-s} \binom{v-s}{t-s} = \sum_{t=s}^v (-1)^{t-s} \binom{v-s}{t-s} - (-1)^{v-s} \binom{v-s}{v-s} = -(-1)^{v-s}$$

then (9) and (10) are equivalent.

The analysis of variance model is obtained from the above results by taking  $v = N$  in equation (10). This gives

$$f(x_1, \dots, x_N) = \bar{g}(x_1, \dots, x_N) - \sum_{t=0}^{N-1} \sum_{(x_{h_1}, \dots, x_{h_t}) \text{ in } (x_1, \dots, x_N)} f(x_{h_1}, \dots, x_{h_t})$$

or, since  $(x_1, \dots, x_N) = x$  and  $\bar{g}(x) = g(x)$ ,

$$g(x) = \sum_{t=0}^N \sum_{(x_{h_1}, \dots, x_{h_t}) \text{ in } x} f(x_{h_1}, \dots, x_{h_t})$$

Hence, the chance variable  $g(X)$  may be expressed in terms of the linear model

$$(11) \quad g(X) = \sum_{t=0}^N \sum_{(h_1, \dots, h_t) \text{ in } (1, \dots, N)} f(X_{h_1}, \dots, X_{h_t}) .$$

Since all effects  $f(X_{h_1}, \dots, X_{h_t})$  are random variables for  $t > 0$  then, by definition, in order to demonstrate that (11) is an analysis of variance model we must show that all effects other than  $\bar{g}$  have zero mean and are uncorrelated. To accomplish this we use the relation (9); thus, the mean value of  $f(X_{i_1}, \dots, X_{i_v})$  is

$$E \left\{ f(X_{i_1}, \dots, X_{i_v}) \right\} = \sum_{t=0}^v (-1)^{v-t} \sum_{(h_1, \dots, h_t) \text{ in } (i_1, \dots, i_v)} E \left\{ \bar{g}(X_{h_1}, \dots, X_{h_t}) \right\} .$$

Since  $E \left\{ \bar{g}(X_{h_1}, \dots, X_{h_t}) \right\} = \bar{g}$  identically and since there are  $\binom{v}{t}$  subsets of  $(i_1, \dots, i_v)$  which contain  $t$  elements, then the right side above becomes

$$\begin{aligned} &= \bar{g} \sum_{t=0}^v (-1)^{v-t} \binom{v}{t} \\ &= \begin{cases} 0 & \text{for } v > 0 \\ \bar{g} & \text{for } v = 0 \end{cases} \end{aligned}$$

To show that for  $(i_1, \dots, i_s) \neq (j_1, \dots, j_t)$ ,  $f(X_{i_1}, \dots, X_{i_s})$  and  $f(X_{j_1}, \dots, X_{j_s})$  are uncorrelated it is then only necessary to show that

$$E \left\{ f(X_{i_1}, \dots, X_{i_s}) \cdot f(X_{j_1}, \dots, X_{j_t}) \right\} = 0 .$$

This we accomplish by first noting that

$$\begin{aligned}
 (12) \quad & E \left\{ f(X_{i_1}, \dots, X_{i_s}) f(X_{j_1}, \dots, X_{j_t}) \right\} \\
 &= E \left[ f(X_{i_1}, \dots, X_{i_s}) \cdot E \left\{ f(X_{j_1}, \dots, X_{j_t}) \mid X_{i_1}, \dots, X_{i_s} \right\} \right] \\
 &= E \left[ f(X_{i_1}, \dots, X_{i_s}) \sum_{q=0}^t (-1)^{t-q} \sum_{(h_1, \dots, h_q) \text{ in } (j_1, \dots, j_t)} E \left\{ \bar{g}(X_{h_1}, \dots, X_{h_q}) \mid X_{i_1}, \dots, X_{i_s} \right\} \right].
 \end{aligned}$$

We next observe that if  $(X_{h_1}, \dots, X_{h_q}) \cap (X_{i_1}, \dots, X_{i_s})$  denotes the (intersection) chance variables which occur in both sets  $(X_{i_1}, \dots, X_{i_s})$

and  $(X_{h_1}, \dots, X_{h_q})$  then, since

$$\begin{aligned}
 E \left\{ \bar{g}(X_{h_1}, \dots, X_{h_q}) \mid X_{i_1}, \dots, X_{i_s} \right\} &= E \left[ E \left\{ g(X) \mid X_{h_1}, \dots, X_{h_q} \right\} \mid X_{i_1}, \dots, X_{i_s} \right] \\
 &= E \left\{ g(X) \mid (X_{h_1}, \dots, X_{h_q}) \cap (X_{i_1}, \dots, X_{i_s}) \right\},
 \end{aligned}$$

the right side of (12) becomes

$$= E \left[ f(X_{i_1}, \dots, X_{i_s}) \sum_{q=0}^t (-1)^{t-q} \sum_{(h_1, \dots, h_q) \text{ in } (j_1, \dots, j_t)} E \left\{ g(X) \mid (X_{h_1}, \dots, X_{h_q}) \cap (X_{i_1}, \dots, X_{i_s}) \right\} \right].$$

Now let  $(v_1, \dots, v_n) = (i_1, \dots, i_s) \cap (j_1, \dots, j_t)$  and consider an arbitrary subset  $(v_{k_1}, \dots, v_{k_m})$  of the set  $(v_1, \dots, v_n)$ . Then among the  $\binom{t}{q}$  subsets of  $(j_1, \dots, j_t)$  which contain  $q$  elements there are  $\binom{t-n}{q-m}$  which contain the set  $(v_{k_1}, \dots, v_{k_m})$  and no other elements of the set  $(v_1, \dots, v_n)$ ,  $m \leq q \leq m + t - n$ . Hence, the coefficient of  $E \left\{ g(X) \mid X_{v_{k_1}}, \dots, X_{v_{k_m}} \right\}$

in the above square bracket is

$$\begin{aligned}
 (13) \quad & f(X_{i_1}, \dots, X_{i_s}) \sum_{q=m}^{m+t-n} (-1)^{t-q} \binom{t-n}{q-m} \\
 & = (-1)^{n-m} f(X_{i_1}, \dots, X_{i_s}) \sum_{q=m}^{m+t-n} (-1)^{(t-n)-(q-m)} \binom{t-n}{q-m}
 \end{aligned}$$

which vanishes when  $t > n$ ; i.e., when the set  $(j_1, \dots, j_t)$  is not a subset of  $(i_1, \dots, i_s)$ . However, since  $(i_1, \dots, i_s) \neq (j_1, \dots, j_t)$  then if  $(j_1, \dots, j_t)$  is a subset of  $(i_1, \dots, i_s)$  it is a proper subset (i.e., not equal to) and the above argument may be applied replacing (12) by

$$\begin{aligned}
 & E \left\{ f(X_{i_1}, \dots, X_{i_s}) f(X_{j_1}, \dots, X_{j_t}) \right\} \\
 & = E \left[ f(X_{j_1}, \dots, X_{j_t}) E \left\{ f(X_{i_1}, \dots, X_{i_s}) \mid X_{j_1}, \dots, X_{j_t} \right\} \right].
 \end{aligned}$$

We may assume, therefore, that  $t > n$ . Since  $(v_{k_1}, \dots, v_{k_m})$  was an arbitrary subset of  $(v_1, \dots, v_n)$  then (13) holds for all subsets of  $(v_1, \dots, v_n)$  so the right side of (12) vanishes completely giving the desired result

$$E \left\{ f(X_{i_1}, \dots, X_{i_s}) f(X_{j_1}, \dots, X_{j_t}) \right\} = 0 \text{ when } (i_1, \dots, i_s) \neq (j_1, \dots, j_t).$$

The variance  $\sigma_g^2$  of the chance variable  $g(X)$  is therefore

$$\begin{aligned}
 (14) \quad \sigma_g^2 &= \sum_{t=1}^N \sum_{(h_1, \dots, h_t) \text{ in } (1, \dots, N)} E \left\{ f(X_{h_1}, \dots, X_{h_t}) \right\}^2.
 \end{aligned}$$

#### Application to the Random Mating Problem

A special case of the above analysis of variance model is obtained by imposing additional restrictions upon the underlying probability model. We have assumed that the chance variables  $X_1, \dots, X_N$  are independently distributed; the additional restrictions which characterize the random mating



problem are, first, that  $N = rn$ , where  $r$  and  $n$  are integers, and second, for every  $i$ ,  $1 \leq i \leq n$ , the chance variables  $X_1, X_{n+1}, \dots, X_{(r-1)n+i}$  are identically distributed. In this special case, then, the  $N$  factors consist of  $r$  identical sets of  $n$  factors. In the genetic problem the  $n$  factors represent  $n$  loci and the  $m_i$  levels,  $A_1^i, \dots, A_{m_i}^i$ , of the  $i$ 'th factor represent the different alleles present at this locus in the population. The integer  $r$ , then, represents the number of repetitions of the haploid chromosomal complement present in the organisms; for  $r > 2$  the organisms must be autopolyploids. We shall consider only the case,  $r = 2$ , of a diploid population at equilibrium under random mating.

Let  $A_{x_1}^1, \dots, A_{x_n}^n$  denote the set of genes contributed by the sire and  $A_{x_{n+1}}^1, \dots, A_{x_{2n}}^n$  denote the set of genes contributed by the dam, so that an individual's genotype is written  $A_{x_1}^1, A_{x_{n+1}}^1, \dots, A_{x_n}^n A_{x_{2n}}^n$ . The genotypic value  $g(x)$  of the genotype defined by  $x = (x_1, \dots, x_{2n})$  may then be expressed as

$$g(x) = \sum_{t=0}^{2n} f(x_{h_1}, \dots, x_{h_t})$$

$(x_{h_1}, \dots, x_{h_t}) \text{ in } x$

and the chance variable  $g(X)$  as

$$g(X) = \sum_{t=0}^{2n} f(X_{h_1}, \dots, X_{h_t})$$

$(h_1, \dots, h_t) \text{ in } (1, \dots, 2n)$

with the understanding that  $f(x_{h_1}, \dots, x_{h_t})|_{t=0} = f(X_{h_1}, \dots, X_{h_t})|_{t=0} = \bar{g}$ .

Since  $X_1$  and  $X_{n+1}$  are identically distributed then for each  $t$ ,  $0 \leq t \leq 2n$ , certain of the chance variables  $f(X_{h_1}, \dots, X_{h_t})$  are also identically distributed. For example, for  $t = 2$  the 4 chance variables  $f(X_1, X_j), f(X_1, X_{n+j}),$

$f(X_{n+i}, X_j)$  and  $f(X_{n+i}, X_{n+j})$ ,  $1 \leq i < j \leq n$ , are identically distributed, and therefore have the same variance; there are  $\binom{n}{2}$  such sets of 4 identically distributed chance variables among the  $\binom{2n}{2}$  chance variables in the sum

$$\sum_{(h_1, h_2) \text{ in } (1, \dots, 2n)} f(X_{h_1}, X_{h_2})$$

and the remaining  $n = \binom{2n}{2} - 4\binom{n}{2}$  chance variables are of the form  $f(X_i, X_{n+i})$ ,  $1 \leq i \leq n$ . We may therefore write

$$\begin{aligned} \sum_{(h_1, h_2) \text{ in } (1, \dots, 2n)} E \left\{ f(X_{h_1}, X_{h_2}) \right\}^2 &= 4 \sum_{(i,j) \text{ in } (1, \dots, n)} E \left\{ f(X_i, X_j) \right\}^2 \\ &+ \sum_{(i) \text{ in } (1, \dots, n)} E \left\{ f(X_i, X_{n+i}) \right\}^2. \end{aligned}$$

The chance effect  $f(X_i, X_j)$ ,  $1 \leq i < j \leq n$ , represents an interaction effect between two genes at different loci as compared to  $f(X_i, X_{n+i})$  which represents an interaction between two alleles at the same locus. Because of their genetic significance these two types of effects have been given different names,  $f(X_i, X_j)$  being called an "additive x additive" effect and  $f(X_i, X_{n+i})$  a "dominance" effect. The expression

$$4 \sum_{(i,j) \text{ in } (1, \dots, n)} E \left\{ f(X_i, X_j) \right\}^2$$

is then called  $\sigma_{A^2}^2$ , the "additive x additive" component of genetic variance, and,

$$\sum_{(i) \text{ in } (1, \dots, n)} E \left\{ f(X_i, X_{n+i}) \right\}^2$$

is called  $\sigma_D^2$ , the "dominance" component of variance. In general, it can be seen that

$$\sum_{(h_1, \dots, h_t) \text{ in } (1, \dots, 2n)} f(X_{h_1}, \dots, X_{h_t})$$

contains  $2^{t-2v}$  chance variables which have the same distribution as

$$f(X_{i_1}, \dots, X_{i_v}, X_{i_v+1}, \dots, X_{i_{t-v}}, X_{n+i_1}, \dots, X_{n+i_v}), 1 \leq i_1 < \dots < i_{t-v} \leq n.$$

Such an effect is called an "additive x ... x additive x dominance x ... x dominance" effect, where the word "additive" appears  $t - 2v$  times and "dominance" appears  $v$  times. Since

$$\begin{aligned} & \sum_{(h_1, \dots, h_t) \text{ in } (1, \dots, 2n)} E \left\{ f(X_{h_1}, \dots, X_{h_t}) \right\}^2 \\ &= \sum_{v=\max(0, t-n)}^{\left[ \frac{t}{2} \right]} 2^{t-2v} \sum_{(i_1, \dots, i_{t-v}) \text{ in } (1, \dots, n)} E \left\{ f(X_{i_1}, \dots, X_{i_v}, \dots, X_{i_{t-v}}, X_{n+i_1}, \dots, X_{n+i_v}) \right\}^2 \end{aligned}$$

where  $\left[ \frac{t}{2} \right]$  = the largest integer in  $\frac{t}{2}$ , then the expression

$$2^{t-2v} \sum_{(i_1, \dots, i_{t-v}) \text{ in } (1, \dots, n)} E \left\{ f(X_{i_1}, \dots, X_{i_v}, \dots, X_{i_{t-v}}, X_{n+i_1}, \dots, X_{n+i_v}) \right\}^2$$

is called  $\sigma_A^{2t-2v} \sigma_D^v$ . Thus, the genetic variance  $\sigma_g^2$  because

$$\sigma_g^2 = \sum_{t=1}^{2n} \sum_{v=\max(0, t-n)}^{\left[ \frac{t}{2} \right]} \sigma_A^{2t-2v} \sigma_D^v$$

or

$$\sigma_g^2 = \sum_{\substack{0 \leq r, s \leq n \\ 1 \leq r + 2s \leq 2n}} \sigma_A^{2r} \sigma_D^{2s}.$$

#### References

Kempthorne, O. (1954). The correlation between relatives in a random mating population. Proc. Roy. Soc. London, B, 143:103-113.

PROBLEM FOR STUDENTS OF PLANT BREEDING 214: SIMPLE AUTOPOLYPLOIDY

If in the analysis of variance model  $N = rn$  and  $X_i, X_{i+n}, \dots, X_{i+(r-1)n}$  are  $r$  identically distributed chance variables,  $i=1, 2, \dots, n$ , then the model may be regarded as describing genotypic values in a simple autopolyploid population which is at equilibrium under random mating. The expression for genetic variance,

$$(1) \quad \sigma_g^2 = \sum_{t=1}^{rn} (h_1, \dots, h_t) \sum_{(1, \dots, rn)} E \left\{ f(X_{h_1}, \dots, X_{h_t}) \right\}^2$$

may again be reduced, as in the diploid case  $r = 2$ , in that for each  $t$  certain of the chance variables  $f(X_{h_1}, \dots, X_{h_t})$  are identically distributed and hence have the same variance. Thus, among the  $\binom{rn}{t}$  chance variables in the sum

$$\sum_{(h_1, \dots, h_t) \text{ in } (1, \dots, rn)} f(X_{h_1}, \dots, X_{h_t})$$

there are  $\prod_{\alpha=1}^r \binom{r}{\alpha}^{k_\alpha}$  chance variables which have the same distribution as

$$f(X_{h_{11}}, \dots, X_{h_{1k_1}}, X_{h_{21}}, X_{h_{21}+n}, \dots, X_{h_{2k_2}}, X_{h_{2k_2}+n}, \dots, X_{h_{r1}}, \dots, X_{h_{r1}+(r-1)n}, \dots, X_{h_{rk_r}}, \dots, X_{h_{rk_r}+(r-1)n})$$

where  $1 \leq h_{\alpha 1} < \dots < h_{\alpha k_\alpha} \leq n$  for  $\alpha = 1, 2, \dots, r$  and  $\sum_{\alpha=1}^r \alpha k_\alpha = t$ . We therefore define the genetic variance component

$$\sigma_{k_1, \dots, k_r}^2 = \prod_{\alpha=1}^r \binom{r}{\alpha}^{k_\alpha} \sum_{(h_{11}, \dots, h_{rk_r}) \text{ in } (1, \dots, n)} E \left\{ f(X_{h_{11}}, \dots, X_{h_{1k_1}}, X_{h_{21}}, X_{h_{21}+n}, \dots, X_{h_{2k_2}}, X_{h_{2k_2}+n}, \dots, X_{h_{r1}}, \dots, X_{h_{r1}+(r-1)n}, \dots, X_{h_{rk_r}}, \dots, X_{h_{rk_r}+(r-1)n}) \right\}^2$$

The problem is to prove that

$$\sigma_g^2 = \sum_{\substack{0 \leq k_1, \dots, k_r \leq n \\ 1 \leq \sum_{\alpha=1}^r \alpha k_\alpha \leq rn}} \sigma_{k_1, \dots, k_r}^2$$

Clearly, by (1), it will suffice to show that for arbitrary  $t$ ,  $1 \leq t \leq rn$ ,

$$(2) \quad \sum_{(h_1, \dots, h_t) \text{ in } (1, \dots, rn)} E \left\{ f(X_{h_1}, \dots, X_{h_t}) \right\}^2 = \sum_{\substack{0 \leq k_1, \dots, k_r \leq n \\ \sum_{\alpha=1}^r \alpha k_\alpha = t}} \sigma_{k_1, \dots, k_r}^2$$

or, in other words, that because of our definition of  $\sigma_{k_1, \dots, k_r}^2$  the right sum of (2) accounts for each of the  $\binom{rn}{t}$  terms appearing in the left sum of (2). The problem thus reduces to one of combinatorial analysis.

Motivation for this problem may be found in the paper "The correlation between relatives in a simple autotetraploid population" by O. Kempthorne, Genetics 40:168-174, 1955. Kempthorne treats the case  $r = 4$ ,  $n = 1$ .

# SOLUTION TO THE AUTOPOLYPLOIDY PROBLEM:

The set  $(h_{11}, \dots, h_{rk_r})$  contains  $\sum_{\alpha=1}^r k_{\alpha}$  integers from the set  $(1, \dots, n)$ . The number of different subsets of  $\sum_{\alpha=1}^r k_{\alpha}$  integers which may be formed from the set  $(1, \dots, n)$  is  $\binom{n}{\sum k_{\alpha}}$ . Associated with each such subset of  $\sum k_{\alpha}$  integers is the collection of chance variables  $f$  which can be formed by partitioning the  $\sum k_{\alpha}$  subscripts  $(h_{11}, \dots, h_{rk_r})$  into  $r$  subsets containing  $k_1, k_2, \dots, k_r$  elements, respectively. The number of such partitions is

$$\frac{(\sum k_{\alpha})!}{\prod (k_{\alpha}!)}$$

The variance component  $\sigma_{k_1, \dots, k_r}^2$  therefore accounts for

$$\prod_{\alpha=1}^r \binom{r}{\alpha}^{k_{\alpha}} \binom{n}{\sum k_{\alpha}} \frac{(\sum k_{\alpha})!}{\prod (k_{\alpha}!)} = \prod_{\alpha=1}^r \binom{r}{\alpha}^{k_{\alpha}} \frac{n!}{\prod (k_{\alpha}!) (n - \sum k_{\alpha})!}$$

of the  $\binom{rn}{t}$  terms in the left sum of (2). It remains, then, only to show that

$$\sum_{\substack{0 \leq k_1, \dots, k_r \leq n \\ \sum k_{\alpha} = t}} \prod_{\alpha=1}^r \binom{r}{\alpha}^{k_{\alpha}} \frac{n!}{\prod (k_{\alpha}!) (n - \sum k_{\alpha})!} = \binom{rn}{t}$$

This we accomplish by comparing coefficients of  $b^t$  on both sides of

$$(1+b)^{rn} = [(1+b)^r]^{\sum k_{\alpha} + n - \sum k_{\alpha}}$$

On the left side we get  $\binom{rn}{t}$  as the coefficient of  $b^t$ . The right side we expand as

$$\begin{aligned} [(1+b)^r]^{\sum k_{\alpha} + n - \sum k_{\alpha}} &= [1 + \binom{r}{1}b + \binom{r}{2}b^2 + \dots + \binom{r}{r}b^r]^{\sum k_{\alpha} + n - \sum k_{\alpha}} \\ &= \sum_{\substack{0 \leq k_1, \dots, k_r \leq n}} \frac{n!}{k_1! \dots k_r! (n - \sum k_{\alpha})!} \prod \binom{r}{\alpha}^{k_{\alpha}} b^{\sum k_{\alpha}} \end{aligned}$$

Hence the coefficient of  $b^{\sum k_{\alpha}} \big|_{\sum k_{\alpha}=t}$  on the right side is

$$\sum_{\substack{0 \leq k_1, \dots, k_r \leq n \\ \sum k_{\alpha} = t}} \frac{n!}{k_1! \dots k_r! (n - \sum k_{\alpha})!} \prod_{\alpha=1}^r \binom{r}{\alpha}^{k_{\alpha}}$$

and the problem is solved.